

# Coarse to fine non-rigid registration: a chain of scale-specific neural networks for multimodal image alignment with application to remote sensing

Armand Zampieri, Guillaume Charpiat, Yuliya Tarabalka

## ► To cite this version:

Armand Zampieri, Guillaume Charpiat, Yuliya Tarabalka. Coarse to fine non-rigid registration: a chain of scale-specific neural networks for multimodal image alignment with application to remote sensing. 2018. hal-01718263

**HAL Id: hal-01718263**

**<https://hal.inria.fr/hal-01718263>**

Preprint submitted on 27 Feb 2018

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Coarse to fine non-rigid registration: a chain of scale-specific neural networks for multimodal image alignment with application to remote sensing

Armand Zampieri<sup>1</sup>, Guillaume Charpiat<sup>2</sup>, Yuliya Tarabalka<sup>1</sup>

<sup>1</sup> *Titane team, INRIA Sophia-Antipolis*

<sup>2</sup> *TAU team, LRI, INRIA Saclay, Université Paris-Sud*

{guillaume.charpiat, yuliya.tarabalka}@inria.fr

*This project started early 2017 and this paper was submitted for publication in November 2017.*

## Abstract

*We tackle here the problem of multimodal image non-rigid registration, which is of prime importance in remote sensing and medical imaging. The difficulties encountered by classical registration approaches include feature design and slow optimization by gradient descent. By analyzing these methods, we note the significance of the notion of scale. We design easy-to-train, fully-convolutional neural networks able to learn scale-specific features. Once chained appropriately, they perform global registration in linear time, getting rid of gradient descent schemes by predicting directly the deformation.*

*We show their performance in terms of quality and speed through various tasks of remote sensing multimodal image alignment. In particular, we are able to register correctly cadastral maps of buildings as well as road polylines onto RGB images, and outperform current keypoint matching methods.*



Figure 1. **Multimodal matching.** We align aerial images with cadastral images. Left: original misregistered images, right: after our realignment.

## 1 Introduction

Image alignment, also named *non-rigid registration*, is the task of finding a correspondence field between two given images, *i.e.* a deformation which, when applied to the first image, warps it to the second one. Such warpings can prove useful in many situations: to transfer information between several images (for instance, from a template image with labeled parts), to compare the appearance of similar parts (as pixel intensity comparison makes sense only after alignment), or to estimate spatial changes (to monitor the evolution of a tumor given a sequence of scans of the same patient over time, for instance). Image alignment has thus been a predominant topic in fields such as medical imaging or remote sensing [30, 19].

### 1.1 Remote sensing & Image alignment

In remote sensing, images of the Earth can be acquired through different types of sensors, in the visible spectrum (RGB) or not (hyperspectral images), from satellites or planes, with various spatial precision (from cm to km range). The analysis of these images allows the monitoring of ecosystems (plants [11], animals [32], soils...) and their evolution (ice melting, drought monitoring, natural disasters and associated help planning), urban growth, as well as the automatic creation of maps [20] or more generally digitizing the Earth.

However, the geographic localization of pixels in remote sensing images is limited by a number of factors, such as the positioning precision and the effect of the relief on non-vertical points of view. The deformation of these images is significant: for instance, in OpenStreetMap [10], objects may be shifted by 8 meters (which is far above the required precision of maps for autonomous driving), which means an error displacement of more than 20 pixels for a 30 cm/pixel resolution.

These deformations prevent a proper exploitation of such data. For instance, let us consider the task of finding buildings and roads in a remote sensing image. While ground truth is actually available in considerable amounts, such as in OpenStreetMap (OSM) based on cadastral information, which gives coordinates (latitude and longitude) of each building corner, this hand-made ground truth is often inaccurate, because of human mistakes. Thus it is not possible to learn from it, as remote sensing images are not properly aligned to it and objects might even not overlap. This is a severe issue for the remote sensing field in the era of big data and machine learning. Many works have been focusing on this problem [2], from the use of relief knowledge to dedicated hand-designed alignment algorithms. Another approach worth mentioning is to train coarsely on the datasets available and fine-tune on small better-hand-aligned datasets [17]. We will here tackle the problem of non-rigid alignment directly.

## 1.2 Classical approaches for non-rigid registration

**Tasks** Image registration deals with images either of the same modality (same sensor), or not. When of the same modality, the task is typically to align different but similar objects (*e.g.*, faces [4] or organs of different people [13]), or to align the same object but taken at different times (as in the tumor monitoring example). On the other side, multi-modal registration deals with aligning images usually of the same object but taken with different sensors, which capture different physical properties, at possibly different resolutions. For instance in medical imaging, MR and CT scans capture the density of water and of matter respectively, while in remote sensing RGB and hyperspectral data capture information from different light frequencies (infrared, etc.). In our case of study, we will focus on the alignment of RGB remote sensing images with cadastral maps, which are vector-format images with polygonal representations of all buildings and roads, hand-made by local authorities, map makers or OpenStreetMap users, as in Figures 1 and 2.

Whether mono-modal or multi-modal, image registration faces two challenges: first, to describe locally image data, and then, to match points with similar description, in a spatially coherent manner. Historically, two main classical approaches have emerged.

**Matching key-points** The first one consists in sampling a few key-points from each image (*e.g.* with Harris corner detection), in describing them locally (with SIFT, SURF, HOG descriptors...) [29, 7], in matching these points [26] and then interpolating to the rest of the image. The ques-

tion is then how to design proper sampling criteria, descriptors, and matching algorithm. In the multi-modal case, one would also have to design or learn the correspondence between the descriptors of the two modalities. Note that high-precision registration requires a dense sampling, as well as consequently finer descriptors, which naturally leads to the second approach.

### Estimating a deformation field by gradient descent

The second approach, particularly popular in medical imaging, consists in estimating a dense deformation field from one image to the other one [1, 9, 25, 15, 13]. One of its advantages over the first approach is to be able to model objects, to make use of shape statistics, etc. The warping is modeled as a smooth vector field  $\phi$ , mapping one image domain onto the other one. Given two images  $I_1$  and  $I_2$ , a criterion  $C(I_1 \circ \phi, I_2)$  is defined, to express the similarity between the warped image  $I_1 \circ \phi$  and the target  $I_2$ , and is optimized with respect to  $\phi$  by gradient descent. Selecting a suitable similarity criterion  $C$  is crucial, as well as designing carefully the gradient descent, as we will detail in section 2.

## 1.3 The new paradigm: neural networks

The difficulty to design or pick particular local descriptors or matching criteria among many possibilities is the trait of computer vision problems where the introduction of neural networks can prove useful. The question is how. Machine learning techniques have already been explored to learn similarity measures between different imaging modalities [38], for instance using kernel methods to register MR and CT brain scans [16], but without tackling the question of scale. We will aim at designing a system able to learn scale-specific and modality-specific features, and able to perform multimodal image registration densely and swiftly, without the use of any iterative process such as gradient descent which hampers classical approaches. Our contributions are thus:

- a swift system to register images densely
- learning features to register images of different modalities
- learning scale-specific features and managing scales
- designing a (relatively small) neural network to do this in an end-to-end fashion
- aligning remote sensing images with cadastral maps (buildings and roads)
- providing a long-awaited tool to create large-scale benchmarks in remote sensing.

We first analyze the problems related to scale when aligning images, in order to design a suitable neural network archi-

texture. We show results on benchmarks and present additional experiments to show the flexibility of the approach.

## 2 Analysis of the gradient descent framework

In order to analyze issues that arise when aligning images, let us first consider the case of mono-modal registration, for simplicity. Keeping the notations from section 1.2, we pursue the quest for a reasonable criterion  $C(I_1 \circ \phi, I_2)$  to optimize by gradient descent to estimate the deformation  $\phi$ .

### 2.1 A basic example

Too local quantities such as the pixelic intensity difference  $C(I_1 \circ \phi, I_2) = \|I_1 \circ \phi - I_2\|_{L^2}^2$  would create many local minima and get the gradient descent stuck very fast. Indeed, if as a toy example one considers  $I_1$  and  $I_2$  to be two white images with a unique black dot<sup>1</sup> at different locations  $\mathbf{x}_1$  and  $\mathbf{x}_2$  respectively, the derivative of  $C(I_1 \circ \phi, I_2)$  with respect to  $\phi$  will never involve quantities based on these two points at the same time, which prevents them from being influenced by each other:

$$\frac{\partial C(I_1 \circ \phi, I_2)}{\partial \phi}(\mathbf{x}) = 2(I_1 \circ \phi(\mathbf{x}) - I_2(\mathbf{x})) (\nabla_{\mathbf{x}} I_1)_{(\phi(\mathbf{x}))}$$

is 0 at all points  $\mathbf{x} \neq \mathbf{x}_1$ , and at  $\mathbf{x}_1$  the deformation  $\phi$  (initialized to the identity) evolves to make it disappear from the cost  $C$  by shrinking the image around. Thus the derivative of the similarity cost  $C$  with respect to the deformation  $\phi$  does not convey any information pushing  $\mathbf{x}_1$  towards  $\mathbf{x}_2$ , but on the contrary will make the descent gradient stuck in this (very poor) local minimum.

Instead of the intensity  $I(\mathbf{x})$ , one might want to consider other local, higher-level features  $L(I)(\mathbf{x})$  such as edge detectors, in order to grasp more meaningful information, and thus minimize a criterion for instance of the form:

$$C(I_1 \circ \phi, I_2) = \|L(I_1 \circ \phi) - L(I_2)\|_{L^2}^2. \quad (1)$$

### 2.2 Neighborhood size

The solution consists in considering local descriptors involving larger spatial neighborhoods, wide enough so that the image domains involved in the computations of  $L(I_1 \circ \phi)(\mathbf{x}_1)$  and  $L(I_2)(\mathbf{x}_2)$  for two points  $\mathbf{x}_1$  and  $\mathbf{x}_2$  to be matched overlap significantly. For instance, the computation of the Canny edge detector is performed over a truncated Gaussian neighborhood, whose size is pre-defined by the standard deviation parameter  $\sigma$ . Another example is the local-cross correlation, which compares the local variations

<sup>1</sup>In the continuous setting, consider a smooth compact-support bump.

of the intensity around  $\mathbf{x}_1$  and  $\mathbf{x}_2$  within a neighborhood of pre-defined size [13]. Another famous example is the mutual information between the histograms of the intensity within a certain window with pre-defined size.

### 2.3 Adapting the scale

In all these cases, the neighborhood size is particularly important: if too small, the gradient descent will get stuck in a poor local minimum, while if too large, the image details might be lost, preventing fine registration. What is actually needed is this neighborhood size to be of the same order of magnitude as the displacement to be found. As this displacement is unknown, the neighborhood size needs to be wide enough during the first gradient steps (possibly covering the full image), and has to decrease with time, for the registration to be able to get finer and finally reach pixelic precision. Controlling the speed of this decrease is a difficult matter, leading to slow optimization. Moreover, the performance of the descriptors may depend on the scale, and different descriptors might need to be chosen for the coarse initial registration than for the finest final one. In plus of the difficult task of designing [39, 36] or learning [16] relevant descriptors  $L_s$  for each scale, this raises another issue, that is that the criterion  $C_s$  to optimize

$$C_s(I_1 \circ \phi, I_2) = \|L_s(I_1 \circ \phi) - L_s(I_2)\|_{L^2}^2 \quad (2)$$

now depends on the current neighborhood size  $s(t)$ , which is itself time-dependent, and thus the optimized criterion  $C_{s(t)}$  might increase when the descriptor  $L_{s(t)}$  evolves: the optimization process is then not a gradient descent anymore.

One might think of scale-invariant descriptors such as SIFT, however the issue is not just to adapt the scale to a particular location within an image, but to adapt it to the amplitude of the deformation that remains to be done to be matched to the *other* image.

### 2.4 Multi-resolution viewpoint

Another point of view on this scale-increasing process is to consider that the descriptors and optimization process remain the same at all scales, but that the *resolution* of the image is increasing. The algorithm is then a loop over successive resolutions [13, 4], starting from a low-resolution version of the image, waiting for convergence of the gradient descent, then upsampling the deformation field found to a higher resolution version of the image, and iterating until the original resolution is reached. The limitation is then that the same descriptor has to be used for each resolution, and, as previously, that the convergence of a gradient descent has to be reached at each scale, leading to slow optimization. A different approach consists in dealing with all

scales simultaneously by considering a multi-scale parameterisation of the deformation [28]. However, the same local minimum problem would be encountered if implemented naively; heuristics then need to be used to estimate at which scale the optimization has currently to be performed locally.

## 2.5 Keeping deformations smooth

Deformations are usually modeled as diffeomorphisms [1, 9, 15, 13], *i.e.* smooth one-to-one vector fields, in order to avoid deleting image parts. The smoothness is controlled by an additional criterion to optimize, quantifying the regularity of the deformation  $\phi$ , such as its Sobolev norm (penalizing fast variations). As in any machine learning technique, this regularity term sets a prior over the space of possible functions (here, deformations), preventing overfitting (here, spatial noise). But once again, the smoothness level required should depend on the scale, *e.g.* prioritizing global translations and rotations at first, while allowing very local moves when converging. This can be handled by suitable metrics on instantaneous deformations [5, 31]; yet in practice these metrics tend to slow down convergence by over-smoothing gradients  $\nabla_{\phi} C$  at finest scales.

## 3 Introducing neural networks

### 3.1 Learning iterative processes

As neural networks have proved useful to replace hand-designed features for various tasks in the literature recently, and convolutional ones (CNN) in particular in computer vision, one could think, for mono-modal image alignment, of training a CNN in the Siamese network setup [3, 6], in order to learn a relevant distance between image patches. The multi-modal version of this would consist in training two CNN (one per modality) with same output size, in computing the Euclidean norm of the difference of their outputs as a dissimilarity measure, and in using that quantity within a standard non-rigid alignment algorithm, such as a gradient descent over (1). For training, this would however require to be able to differentiate the result of this iterative alignment process with respect to the features. This is not realistic, given the varying, usually large number of steps required for typical alignment tasks. A similar approach was nonetheless successfully used in [17], for the simpler task of correcting blurry segmentation maps, sharpening them and relying on image edges. For this, a partial differential equation (PDE) was mimicked with a recurrent network, and the number of steps applying this PDE was pre-defined to a small value (5), sufficient for that particular problem. In the same spirit, for image denoising, in [21, 34] the proximal operator used during an iterative optimization process is modeled by a neural network and learned. There is to our

knowledge no neural network-based work for dense non-rigid alignment yet. In [24], the Siamese network idea is used, but for matching only very few points. It is also worth noting that, much earlier, in [16], a similarity criterion between different modalities was learned, with a kernel method, but for rigid registration only.

### 3.2 A more direct approach

As seen in the previous sections, aligning images thanks to a gradient descent over the deformation  $\phi$  has the following drawbacks: it is slow because of the need to ensure convergence at each scale, it is actually not a real gradient descent if descriptors are scale-dependent, and it induces a long backpropagation chain when learning the descriptors. To get rid of this iterative process, we propose to predict directly its final result at convergence. That is, given images  $I_1$  and  $I_2$ , to predict directly the optimal deformation  $\phi$  so that  $I_1 \circ \phi$  and  $I_2$  are aligned. Also, instead of proceeding in two steps: first learning the features  $L$  required to define the criterion  $C$  in (1), then finding the deformation  $\phi$  minimizing  $C$ , we propose to directly learn the deformation as in a standard machine learning setup, that is, from examples. Given a training dataset of input pairs  $P = (I_1, I_2)$  together with the expected associated output  $\phi_P$ , we aim at learning the function  $P \mapsto \phi_P$ .

### 3.3 Machine learning setting

**Training set** We first consider the task of aligning geolocalized aerial RGB images with binary maps from OpenStreetMap indicating building locations. As explained in section 1.1, the matching is usually imperfect. Creating the deformation ground truth by manually performing the warpings would be too time-consuming. Instead, we extract image pairs which visually look already well aligned, as in Figure 2. This way we obtain a dataset composed of 5000x5000 image pairs (aerial RGB image, binary vector-format building map) at resolution 0.3m/pixel, for which the deformation  $\phi$  to be found is the identity.

We generate an artificial training set by applying random deformations to the cadastral vectorial maps, moving accordingly the corners of the polygons it contains, and then generating the new binary maps by rasterization. We thus obtain a training set of pairs of non-registered images, with known deformations. As typical deformations in reality are smooth, we model our family of random deformations as: a global translation  $\mathbf{v}_0$  taken uniformly within a certain range  $[-r, +r]^2$ , plus a mixture of Gaussian functions with random shifts  $\mathbf{v}_i$ , centers  $\mathbf{x}_i$  and covariance matrices  $S_i$ :

$$\phi(\mathbf{x}) = \mathbf{v}_0 + \sum_{i=1}^n \mathbf{v}_i e^{-(\mathbf{x}-\mathbf{x}_i) S_i (\mathbf{x}-\mathbf{x}_i)} \quad (3)$$





Figure 2. **Multimodal pair** of already satisfyingly aligned images, from the database. Left: aerial RGB image, right: binary vector-format cadastral image (buildings are shown in white).

with uniformly random  $\mathbf{v}_i$ ,  $S_i$ ,  $\mathbf{x}_i$  within suitable predefined ranges ( $S_i$  being symmetric positive definite). This way, we can drastically augment the dataset by applying arbitrarily many random deformations to the initially well-aligned images.

**Optimization criterion** The loss considered is simply the Euclidean norm of the prediction error:

$$C(\mathbf{w}) = \mathbb{E}_{(I_1, I_2, \phi_{\text{GT}}) \in \mathcal{D}} \left[ \sum_{\mathbf{x} \in \Omega(I_2)} \left\| \hat{\phi}_{(\mathbf{w})(I_1, I_2)}(\mathbf{x}) - \phi_{\text{GT}}(\mathbf{x}) \right\|_2^2 \right]$$

*i.e.* the expectation, over the ground truth dataset  $\mathcal{D}$  of triplet examples (RGB image  $I_1$ , cadastral image  $I_2$ , associated deformation  $\phi_{\text{GT}}$ ), of the sum, over all pixels  $\mathbf{x}$  in the image domain  $\Omega(I_2)$ , of the norm of the difference between the ground truth deformation  $\phi_{\text{GT}}(\mathbf{x})$  and the one predicted  $\hat{\phi}_{(\mathbf{w})(I_1, I_2)}(\mathbf{x})$  for the pair of images  $(I_1, I_2)$  given model parameters  $\mathbf{w}$  (*i.e.* the neural network weights).

In order to make sure that predictions are smooth, we also consider for each pixel a penalty over the norm of the (spatial) Laplacian of the deformation  $\hat{\phi}$ :

$$+ \left\| \Delta \hat{\phi}_{(\mathbf{w})(I_1, I_2)}(\mathbf{x}) \right\|_2^2 \quad (4)$$

which penalizes all but affine warpings. In practice in the discrete setting this sum is the deviation of  $\hat{\phi}(\mathbf{x})$  from the average over the 4 neighboring pixels:  $+ \left\| \hat{\phi}(\mathbf{x}) - \frac{1}{4} \sum_{\mathbf{x}' \sim \mathbf{x}} \hat{\phi}(\mathbf{x}') \right\|_2^2$ .

### 3.4 A first try

We first produce a training set typical of real deformations by picking a realistic range  $r = \pm 20$  pixels of deformation amplitudes. We consider a fully-convolutional neural network, consisting of two convolutional networks (one for each input image  $I_i$ ), whose final outputs are concatenated and sent to more convolutional layers. The last layer has

two features, *i.e.* emits two real values per pixel, which are interpreted as  $\hat{\phi}(\mathbf{x})$ . In our experiments, such a network does not succeed in learning deformations: it constantly outputs  $\hat{\phi}(\mathbf{x}) = (0, 0) \forall \mathbf{x}$ , which is the best constant value for our loss, *i.e.* the best answer one can make when not understanding the link between the input  $(I_1, I_2)$  and the output  $\phi$  for a quadratic loss: the average expected answer  $\mathbb{E}_{(I_1, I_2, \phi_{\text{GT}}) \in \mathcal{D}} [\phi]$ , which is  $(0, 0)$  in our case.

We also tried changing representation by predicting bin probabilities  $p(\Phi_x(\mathbf{x}) \in [a, a+1])$ ,  $p(\Phi_y(\mathbf{x}) \in [b, b+1])$  for each integer  $-r \leq a, b < r$ , by outputting 2 vectors of  $2r$  real values per pixel, but this lead to the same result.

### 3.5 Dealing with a single scale

The task in section 3.4 is indeed too hard: the network needs to develop local descriptors at all scales to capture all information, and is asked to perform a fine matching with  $(2r)^2 \simeq 1700$  possibilities for each pixel  $\mathbf{x}$ .

This task can be drastically simplified though, by requiring the network to perform the alignment *at one scale only*. By this, we mean:

**Task at scale  $s$ :** Solve the alignment problem for the image pair  $(I_1, I_2)$ , with a precision required of  $\pm 2^s$  pixels, under the assumption that the amplitude of the registration to be found is not larger than  $2^{s+1}$  pixels.

For instance, at scale  $s = 0$ , the task is to search for a pixelwise precise registration ( $\pm 1$  pixel) on a dataset prepared as previously but with amplitude  $r = 2^{s+1} = 2$ . As a first approximate test, we train the same network as described earlier, with  $r = 2$ , *i.e.* each of the 2 coordinates of  $\phi(\mathbf{x})$  take value in  $[-2, 2]$ , and we consider a prediction  $\hat{\phi}(\mathbf{x})$  to be correct if in the same unit-sized bin as  $\phi(\mathbf{x})$ . Without tuning the architecture or the optimization method, we obtain, from the first training run, about 90% of accuracy, to be compared to  $\sim 6\%$  for a random guess.

Thus, it is feasible, and easy, to extract information when specifying the scale. Intuitively, the search space is much smaller; in the light of section 2.2, the descriptor receptive field required for such a  $\pm 1$  pixel task is just of radius 1. And indeed, in the classical framework for monomodal registration, a feature as simple as the image intensity would define a sufficiently good criterion (1), as the associated gradient step involves the comparison to the next pixel (through  $\nabla_{\mathbf{x}} I_1$ ). Note that such a simple intensity-based criterion would not be expected to perform more, *e.g.* find deformations of amplitude  $r \geq 2$  pixels in the general case (textures).

**Designing a suitable neural network architecture** We now propose better architectures to solve that alignment

task at scale  $s = 0$ . We need a fully-convolutional architecture since the output is a 2-channel image of the same size as the input, and we need to go across several scales in order to understand which kind of object part each pixel belongs to, in each modality. High-level features require a wide receptive field, and are usually obtained after several pooling layers, in a pyramidally shaped network. The output needs however to be of the same resolution as the input, which leads to autoencoder-like shapes. In order to keep all low-level information until the end, and not to lose precision, we link same-resolution layers together, thus obtaining a U-net-like network [33] (developed for medical image segmentation). As the 2 input images are not registered, and to get modality-specific features, we build 2 separate convolutional pyramids, one for each modality (in a similar fashion as networks for stereo matching [40]), but concatenate their activities once per scale to feed the U-net. The architecture is summarized in Figure 3. The network is trained successfully to solve the  $s = 0$  task as explained previously.

### 3.6 A chain of scale-specific neural networks

We now solve the general alignment task very simply:

**Solution for task at scale  $s$ :** *Downsample the images by a factor  $2^s$ ; solve the alignment task at scale 0 for these reduced images, and upsample the result with the same factor.*

**Full alignment algorithm:** *Given an image pair  $(I_1, I_2)$  of width  $w$ , iteratively solve the alignment task at scale  $s$ , from  $s = \log_2 w$  until  $s = 0$ .*

One can choose to use the same network for all scales, or different ones if we expect specific features at each scale, as in remote sensing or medical imaging.

The full processing chain is shown in Figure 4. Note a certain global similarity with ResNet [12], in that we have a chain of consecutive scale-specific blocks, each of which refining the previously-estimated deformation, not by adding to it but by diffeomorphism composition:  $\phi_{s-1} = \phi_s \circ (\text{Id} + f(I_1 \circ \phi_s, I_2 \circ \phi_s))$ . Another difference with ResNet is that we train each scale-specific block independently, which is much easier than training the whole chain at once.

Note that the overall complexity of an alignment is very low, linear in the image size. Indeed, for a given image with  $n$  pixels, a similar convolutional architecture is applied to all reduced versions by factors  $2^s$ , of size  $2^{-s} \times 2^{-s} n$  pixels, leading to a total cost of  $n(1 + \frac{1}{4} + \frac{1}{16} + \frac{1}{64} + \dots)K < \frac{4}{3}nK$  where  $K$  is the constant per-pixel convolutional cost. This is to be compared with the classical gradient descent based approaches of unknown convergence duration, and with

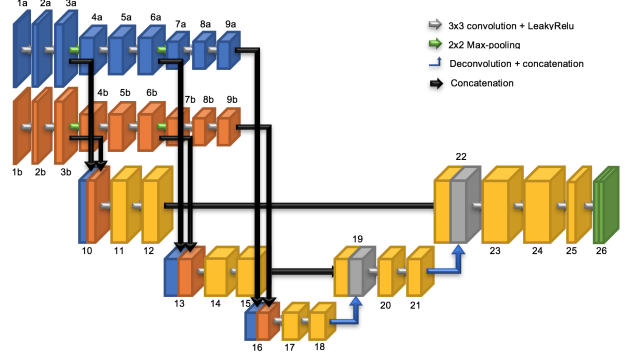


Figure 3. **Network architecture.** The two input images  $I_1$  and  $I_2$  are fed to layers 1a and 1b respectively. The output is a 2 dimensional vector map (layer 26 with 2 channels). See Appendix for all details.

the classical multi-resolution approaches with gradient descents at each scale.

Note also some similarity with recent work on optical flow [14], consisting in an arrangement of 3 different scale-related blocks, though monomodal, not principled from a scale analysis and without scale-specific training.

We will also check the following variations:

- “fast”: replace all scale-specific blocks with the same  $s = 2$ -specific block, to see how well features generalize across scales; the output quality decreases slightly but remains honorable.
- “accurate”: apply the network on symmetrised and rotated versions of the input images, and average the result over these 8 tests. This ensures rotation invariance and improves the result.

## 4 Experiments

We perform four experiments on different datasets. The **first experiment** uses the **Inria aerial image labeling dataset** [18], which is a collection of aerial orthorectified color (RGB) imagery with a spatial resolution of 30 cm/pixel covering 810 km<sup>2</sup> over 9 cities in USA and Austria. We aim at aligning a map of buildings downloaded from OSM with the images from the Inria dataset. The network described in Section 3.6 is trained using image patches from six different cities, for which accurate building cadastral data are available<sup>2</sup>. We then evaluated the network by using images of the area of Kitsap County not presented during training. Fig. 5 shows an example close-up of alignment result.

<sup>2</sup>The cadastral data are extracted from OSM and contain a small misalignment of an order of several pixels.

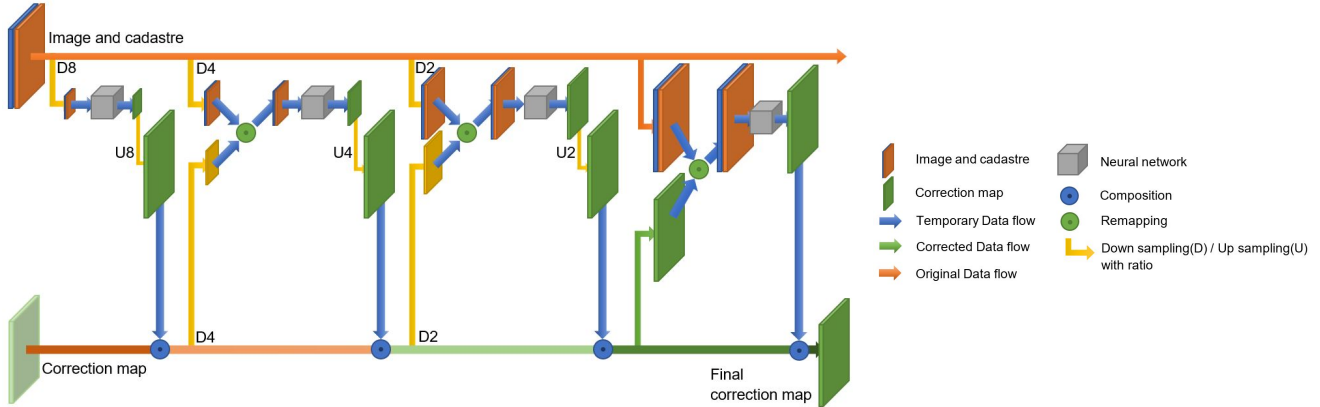


Figure 4. **Full architecture as a chain of scale-specific neural networks.** The two full-resolution input images are always available on the top horizontal row. The full-resolution deformation is iteratively refined, scale per scale, on the bottom horizontal line. Each scale-specific block downsample the images to the right size, apply the previously-estimated deformation, and refines it, in a way somehow similar to ResNet. Each block is trained independently.

In the **second experiment**, the network trained in the first experiment is used to align the OSM building map with satellite images with a pansharpened resolution of 50 cm/pixel, acquired by the Pléiades sensor over the Forez rural area in France.

To measure performance of the network, we use the percentage of correct key point metric [27]. We manually identified matching key points on two couples of multimodal images (one Kitsap image from experiment 1 and one Forez image from experiment 2) with more than 600 keypoints for each image. We then measure the distance between the positions of keypoints after alignment by using different algorithms and the manually indicated ones. If this distance is smaller than a certain threshold, the keypoint is identified as matched.

Fig. 6 compares the performance of our network with the following methods: DeepFlow of Weinzaepfel *et al.* [35], two variations of geometric matching method of Rocco *et al.* [27], and a multimodal registration method of Ye *et al.* [37]. Our approach clearly outperforms other ones by a large margin. We note that averaging over rotations and symmetries (in green, “*accurate*”) does help on the Forez dataset, and that learning scale-specific features performs slightly better than scale-independent features but not always (blue vs. red, “*fast*”). Examples of alignment results are shown in Figure 7. Our approach is also much faster, as shown by the computational times below for a  $5000 \times 5000$  image, even though we compute a dense registration while other approaches only match keypoints:

| Method | Ours         | [27]         | [35]        | [37]   |
|--------|--------------|--------------|-------------|--------|
| Time   | <b>80 s</b>  | 238 s        | 784 s       | 9550 s |
| CPU    | Opteron 2Ghz | Intel 2.7Ghz | Int. 3.5Ghz |        |
| GPU    | GTX 1080 Ti  | Q.M2000M     | GT 960 M    |        |

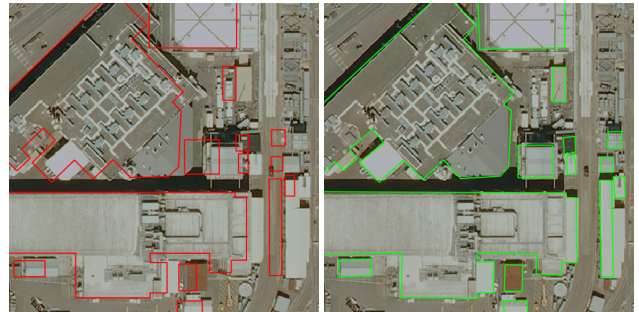
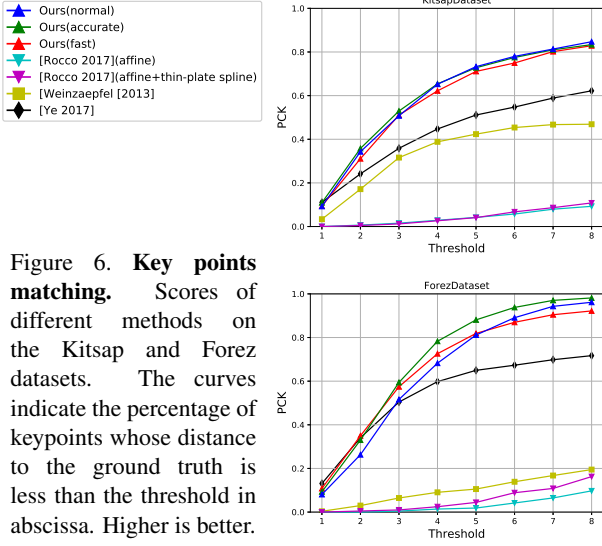


Figure 5. **Example of image alignment.** Left : original image and OpenStreetMap (OSM) map. Right : Alignment result.

In a **third experiment**, we align roads with the images used in the first experiment. The task differs from previous experiments in that only the center line of the road is known, and in the form of a polyline. Moreover, local edges are not useful features for alignment anymore, as the center of roads is homogeneous. We train on OSM data, by dilating road polylines to reach a 4 pixel width and rasterizing them. We then test the performance of the trained network on the Kitsap image. The results are shown in Figure 8.

The **fourth experiment** checks the performance of our approach on a higher-resolution dataset. We consider the Kitti dataset [8], which contains high precision aerial images (9 cm/pixel), as well as perfectly aligned multi classes labeling [20]. We create a training set with artificial random deformations, in the same spirit as before, and a test set with randomly deformed images as well, but following different distributions, in order to check also the robustness of our training approach. Image pairs to be registered consist of a RGB image and a 3-channel binary image indicating build-





ings, roads and sidewalk presence respectively. An example of result is shown in Figure 9. We also analyse the distribution of misalignments before and after registration, shown as histograms in Figure 10. We note that the vast majority of pixels are successfully very closely matched to their ground truth location.

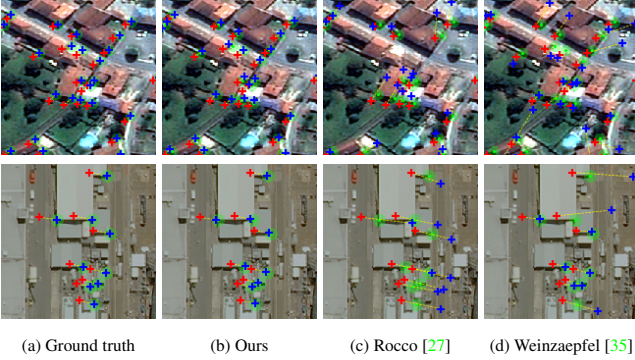


Figure 7. **Multimodal keypoint matching comparison** for different methods and two datasets. Top: Forez dataset; bottom: Kitsap. Blue: predicted, green: ground truth. Full resolution in Appendix.

We also perform an extra experiment to show that our multi-scale approach could generalize to other applications. We consider the problem of stereovision, where the input is a pair of RGB images taken from slightly different view points, and the expected output is the depth map, *i.e.* a single channel image instead of a deformation field. We consider the dataset from [22, 23] and define the loss function as the depth error (squared), plus the regularizer (4). We keep the same architecture but link the scale-specific networks with additions instead of compositions, so that each block adds



Figure 8. **Example of road alignment.** Left: original alignment between image and roads (Kitsap); right: results after realignment.



Figure 9. **Example of alignment on the Kitti Dataset.** Left: before alignment; right: after alignment.

scale-specific details to the depth map. The promising result (first run, no parameter tuning) is shown in the Appendix.

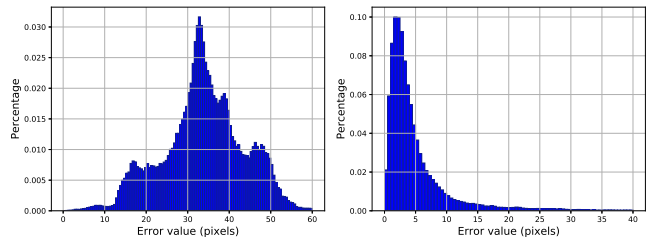


Figure 10. **Misalignment histograms on the fourth experiment (Kitti dataset).** Left: original misalignment distribution in the test set; right: remaining error after our automatic alignment.

**Optimization details** The network is trained with an Adam optimizer, on mini-batches of 16 patches of  $128 \times 128$  pixels images, with a learning rate starting from 0.001

and decayed by 4% every 1000 iterations. Weights are initialized following Xavier Glorot’s method. We trained for 60 000 iterations. More technical details are available in the Appendix.

**Additional details specific to sparse modalities** such as cadastral maps, though not essential. During training, we sort out fully or mostly blank images (*e.g.* cadastre without any building). Also, in order to train more where there is more information to extract (*e.g.*, corners and edges vs. wide homogeneous spaces), we multiply the pixel loss by a factor  $> 1$  on building edges when training.

When rectangular building are glued together in a row with shared walls, the location of their edges and corners is not visible anymore on the rasterized version of the OSM cadastre. By adding a channel to the cadastre map, reminding the OSM corner locations, we observe a better alignment of such rows.

## 5 Conclusion

Based on an analysis of classical methods, we designed a chain of scale-specific neural networks for non-rigid image registration. By predicting directly the final registration at each scale, we avoid slow iterative processes such as gradient descent schemes. The computational complexity is linear in the image size, and far lower than even keypoint matching approaches. We demonstrated its performance on various remote sensing tasks and resolutions. The trained network as well as the training code will be made available online. This way, we hope to contribute to the creation of large datasets in remote sensing, where precision so far was an issue requiring hand-made ground truth.

An interesting point to study further is the specialization of the networks for each scale, to check on which type of image dataset it is strong or not (medical imaging vs. landscape pictures *e.g.*).

## References

- [1] M. F. Beg, M. I. Miller, A. Trouvé, and L. Younes. Computing large deformation metric mappings via geodesic flows of diffeomorphisms. *International journal of computer vision*, 61(2):139–157, 2005. 2, 4
- [2] B. Bischke, P. Helber, J. Folz, D. Borth, and A. Dengel. Multi-task learning for segmentation of building footprints with deep neural networks. *arXiv preprint arXiv:1709.05932*, 2017. 2
- [3] J. Bromley, I. Guyon, Y. LeCun, E. Säckinger, and R. Shah. Signature verification using a” siamese” time delay neural network. In *Advances in Neural Information Processing Systems*, pages 737–744, 1994. 4
- [4] G. Charpiat, R. Keriven, and O. Faugeras. Image statistics based on diffeomorphic matching. In *10th International Conference on Computer Vision*, volume 1, pages 852–857, 2005. 2, 3
- [5] G. Charpiat, P. Maurel, J.-P. Pons, R. Keriven, and O. Faugeras. Generalized gradients: Priors on minimization flows. *International Journal of Computer Vision*, 2007. 4
- [6] S. Chopra, R. Hadsell, and Y. LeCun. Learning a similarity metric discriminatively, with application to face verification. In *Computer Vision and Pattern Recognition, 2005. CVPR 2005. IEEE Computer Society Conference on*, volume 1, pages 539–546. IEEE, 2005. 4
- [7] N. Dalal and B. Triggs. Histograms of oriented gradients for human detection. In *2005 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR’05)*, volume 1, pages 886–893 vol. 1, June 2005. 2
- [8] A. Geiger, P. Lenz, C. Stiller, and R. Urtasun. Vision meets robotics: The kitti dataset. *The International Journal of Robotics Research*, 32(11):1231–1237, 2013. 7
- [9] J. Glaunes, A. Trouvé, and L. Younes. Diffeomorphic matching of distributions: A new approach for unlabelled point-sets and sub-manifolds matching. In *Computer Vision and Pattern Recognition, 2004. CVPR 2004. Proceedings of the 2004 IEEE Computer Society Conference on*, volume 2, pages II–II. Ieee. 2, 4
- [10] M. Haklay and P. Weber. Openstreetmap: User-generated street maps. *IEEE Pervasive Computing*, 7(4):12–18, Oct 2008. 1
- [11] M. C. Hansen, P. V. Potapov, R. Moore, M. Hancher, S. A. Turubanova, A. Tyukavina, D. Thau, S. V. Stehman, S. J. Goetz, T. R. Loveland, A. Kommareddy, A. Egorov, L. Chini, C. O. Justice, and J. R. G. Townshend. High-resolution global maps of 21st-century forest cover change. *Science*, 342(6160):850–853, 2013. 1
- [12] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. *CoRR*, abs/1512.03385, 2015. 6
- [13] G. Hermosillo, C. Chéfd’Hotel, and O. Faugeras. Variational methods for multimodal image matching. *International Journal of Computer Vision*, 50(3):329–343, 2002. 2, 3, 4
- [14] E. Ilg, N. Mayer, T. Saikia, M. Keuper, A. Dosovitskiy, and T. Brox. FlowNet 2.0: Evolution of optical flow estimation with deep networks. *arXiv preprint arXiv:1612.01925*, 2016. 6
- [15] D. G. Kendall. A survey of the statistical theory of shape. *Statistical Science*, pages 87–99, 1989. 2, 4
- [16] D. Lee, M. Hofmann, F. Steinke, Y. Altun, N. D. Cahill, and B. Scholkopf. Learning similarity measure for multimodal 3d image registration. In *Computer Vision and Pattern Recognition, 2009. CVPR 2009. IEEE Conference on*, pages 186–193. IEEE, 2009. 2, 3, 4
- [17] E. Maggiori, G. Charpiat, Y. Tarabalka, and P. Alliez. Recurrent neural networks to enhance satellite image classification maps. *CoRR*, abs/1608.03440, 2016. 2, 4
- [18] E. Maggiori, Y. Tarabalka, G. Charpiat, and P. Alliez. Can semantic labeling methods generalize to any city? the inria aerial image labeling benchmark. In *IEEE International*

*Geoscience and Remote Sensing Symposium (IGARSS)*. IEEE, 2017. 6

- [19] J. B. A. Maintz, P. A. van den Elsen, and M. A. Viergever. Evaluation of ridge seeking operators for multimodality medical image matching. *IEEE Transactions on pattern analysis and machine intelligence*, 18(4):353–365, 1996. 1
- [20] G. Mátyus, S. Wang, S. Fidler, and R. Urtasun. Hd maps: Fine-grained road segmentation by parsing ground and aerial images. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3611–3619, 2016. 1, 7
- [21] T. Meinhardt, M. Möller, C. Hazirbas, and D. Cremers. Learning proximal operators: Using denoising networks for regularizing inverse imaging problems. In *International Conference on Computer Vision*, October 2017. 4
- [22] M. Menze and A. Geiger. Object scene flow for autonomous vehicles. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015. 8
- [23] M. Menze, C. Heipke, and A. Geiger. Joint 3d estimation of vehicles and scene flow. In *ISPRS Workshop on Image Sequence Analysis (ISA)*, 2015. 8
- [24] N. Merkle, W. Luo, S. Auer, R. Miller, and R. Urtasun. Exploiting deep matching and sar data for the geo-localization accuracy improvement of optical satellite images. *Remote Sensing*, 9(6), 2017. 4
- [25] P. W. Michor, D. Mumford, J. Shah, and L. Younes. A metric on shape space with explicit geodesics. *arXiv preprint arXiv:0706.4299*, 2007. 2
- [26] K. Mikolajczyk and C. Schmid. Scale & affine invariant interest point detectors. *International journal of computer vision*, 60(1):63–86, 2004. 2
- [27] I. Rocco, R. Arandjelovic, and J. Sivic. Convolutional neural network architecture for geometric matching. *CoRR*, abs/1703.05593, 2017. 7, 8
- [28] J. A. Schnabel, D. Rueckert, M. Quist, J. M. Blackall, A. D. Castellano-Smith, T. Hartkens, G. P. Penney, W. A. Hall, H. Liu, C. L. Truitt, et al. A generic framework for non-rigid registration based on non-uniform multi-level free-form deformations. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 573–581. Springer, 2001. 4
- [29] P. Scovanner, S. Ali, and M. Shah. A 3-dimensional sift descriptor and its application to action recognition. In *Proceedings of the 15th ACM International Conference on Multimedia*, MM ’07, pages 357–360, New York, NY, USA, 2007. ACM. 2
- [30] A. Sotiras, C. Davatzikos, and N. Paragios. Deformable medical image registration: A survey. *IEEE transactions on medical imaging*, 32(7):1153–1190, 2013. 1
- [31] G. Sundaramoorthi, A. Yezzi, and A. Mennucci. Coarse-to-fine segmentation and tracking using sobolev active contours. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 30(5):851–864, 2008. 4
- [32] Y. Verdie and F. Lafarge. Efficient Monte Carlo sampler for detecting parametric objects in large scenes. In *Proc. of the European Conference on Computer Vision (ECCV)*, Firenze, Italy, 2012. 1
- [33] T. Von Eicken, A. Basu, V. Buch, and W. Vogels. U-net: A user-level network interface for parallel and distributed computing. In *ACM SIGOPS Operating Systems Review*, volume 29, pages 40–53. ACM, 1995. 6
- [34] S. Wang, S. Fidler, and R. Urtasun. Proximal deep structured models. In *Advances in Neural Information Processing Systems*, pages 865–873, 2016. 4
- [35] P. Weinzaepfel, J. Revaud, Z. Harchaoui, and C. Schmid. DeepFlow: Large displacement optical flow with deep matching. In *IEEE International Conference on Computer Vision (ICCV)*, Sydney, Australia, Dec. 2013. 7, 8
- [36] Y. Ye and J. Shan. A local descriptor based registration method for multispectral remote sensing images with non-linear intensity differences. *ISPRS Journal of Photogrammetry and Remote Sensing*, 90:83–95, 2014. 3
- [37] Y. Ye, J. Shan, L. Bruzzone, and L. Shen. Robust registration of multimodal remote sensing images based on structural similarity. *IEEE Transactions on Geoscience and Remote Sensing*, 55(5):2941–2958, 2017. 7
- [38] Y. Ye and L. Shen. Hopc: A novel similarity metric based on geometric structural properties for multi-modal remote sensing image matching. In *Proc. Ann. Photogram., Remote Sens. Spatial Inf. Sci.(ISPRS)*, pages 9–16, 2016. 2
- [39] L. Yu, D. Zhang, and E.-J. Holden. A fast and fully automatic registration approach based on point features for multi-source remote-sensing images. *Computers & Geosciences*, 34(7):838–848, 2008. 3
- [40] J. Zbontar and Y. LeCun. Stereo matching by training a convolutional neural network to compare image patches. *CoRR*, abs/1510.05970, 2015. 6

## A Appendix

### A.1 Neural network architecture details

See Figures 11 and 12 for further details about the architecture and meta-parameters.

### A.2 Alignment framework

The whole processing framework for the alignment of OpenStreetMap cadastral information with aerial images is summarized in the chart shown in Figure 13.

### A.3 Example of deformation

As explained in the article, to augment the dataset size for training, we generate random deformations, as a mixture of Gaussian functions with random shifts. An example of such a deformation (amplified 4 times for better visualisation purposes), and the result of its application to the original image, are shown in Figure 15.

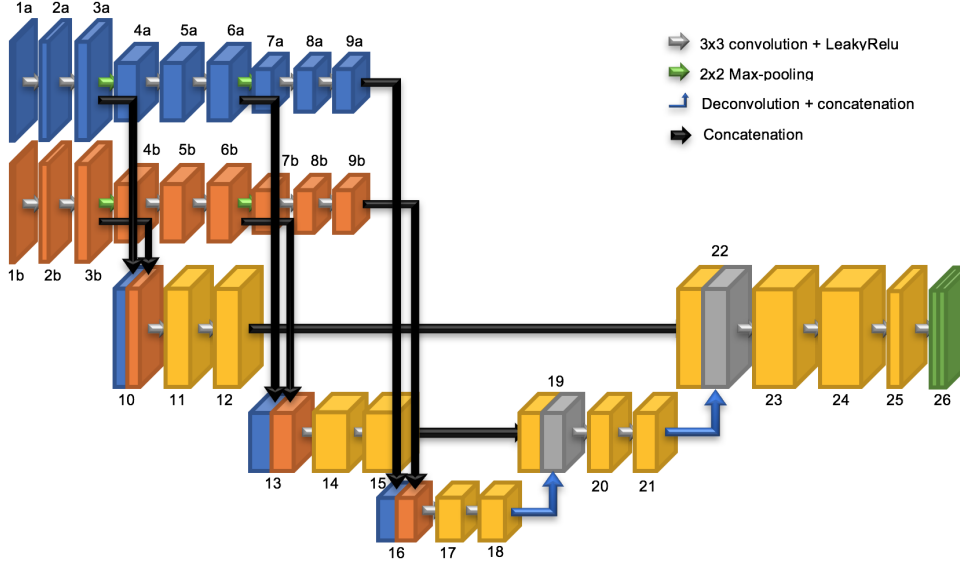


Figure 11. **Network architecture.** The two input images  $I_1$  and  $I_2$  are fed to layers 1a and 1b respectively. The output is a 2 dimensional vector map (layer 26 with 2 channels). This architecture allows to merge information from both sources at all scales, to extract high-level information, and to remember fine details from the input resolution to output a precise full-resolution deformation. Details on Figure 12.

| Start layer | End Layer | Name            | Kernel size | Number of filters | padding | stride |
|-------------|-----------|-----------------|-------------|-------------------|---------|--------|
| 1           | 2         | convolution-1   | 5           | 16                | 2       | 2      |
| 2           | 3         | convolution-2   | 5           | 32                | 2       |        |
| 3           | 4         | pooling-1       | 2           |                   |         |        |
| 4           | 5         | convolution-3   | 3           | 32                | 1       | 2      |
| 5           | 6         | convolution-4   | 3           | 32                | 1       |        |
| 6           | 7         | pooling-2       | 2           |                   |         |        |
| 7           | 8         | convolution-5   | 3           | 64                | 1       |        |
| 8           | 9         | convolution-6   | 3           | 64                | 1       |        |
| 3           | 10        | concatenation-1 |             |                   |         |        |
| 6           | 13        | concatenation-2 |             |                   |         |        |
| 9           | 16        | concatenation-3 |             |                   |         |        |
| 10          | 11        | convolution-7   | 3           | 32                | 1       |        |
| 11          | 12        | convolution-8   | 3           | 32                | 1       |        |
| 13          | 14        | convolution-9   | 3           | 64                | 1       |        |
| 14          | 15        | convolution-10  | 3           | 64                | 1       |        |
| 16          | 17        | convolution-11  | 3           | 64                | 1       |        |
| 18          | 19        | convolution-12  | 3           | 64                | 1       |        |
| 18          | 18'       | deconvolution-1 | 3           | 64                |         |        |
| 15-18'      | 19        | concatenation-4 |             |                   |         |        |
| 19          | 20        | convolution-11  | 3           | 64                | 1       |        |
| 20          | 21        | convolution-12  | 3           | 64                | 1       |        |
| 21          | 21'       | deconvolution-2 | 3           | 32                |         |        |
| 12-21'      | 22        | concatenation-5 |             |                   |         |        |
| 22          | 23        | convolution-13  | 3           | 64                | 1       |        |
| 23          | 24        | convolution-14  | 3           | 64                | 1       |        |
| 24          | 25        | convolution-15  | 3           | 32                | 1       |        |
| 25          | 26        | convolution-16  | 3           | 2                 | 1       |        |

Figure 12. Details for each layer of the (scale-specific) neural network displayed on Figure 3. “Kernel size 3” for a convolutional layer means “ $3 \times 3$ ” convolution.



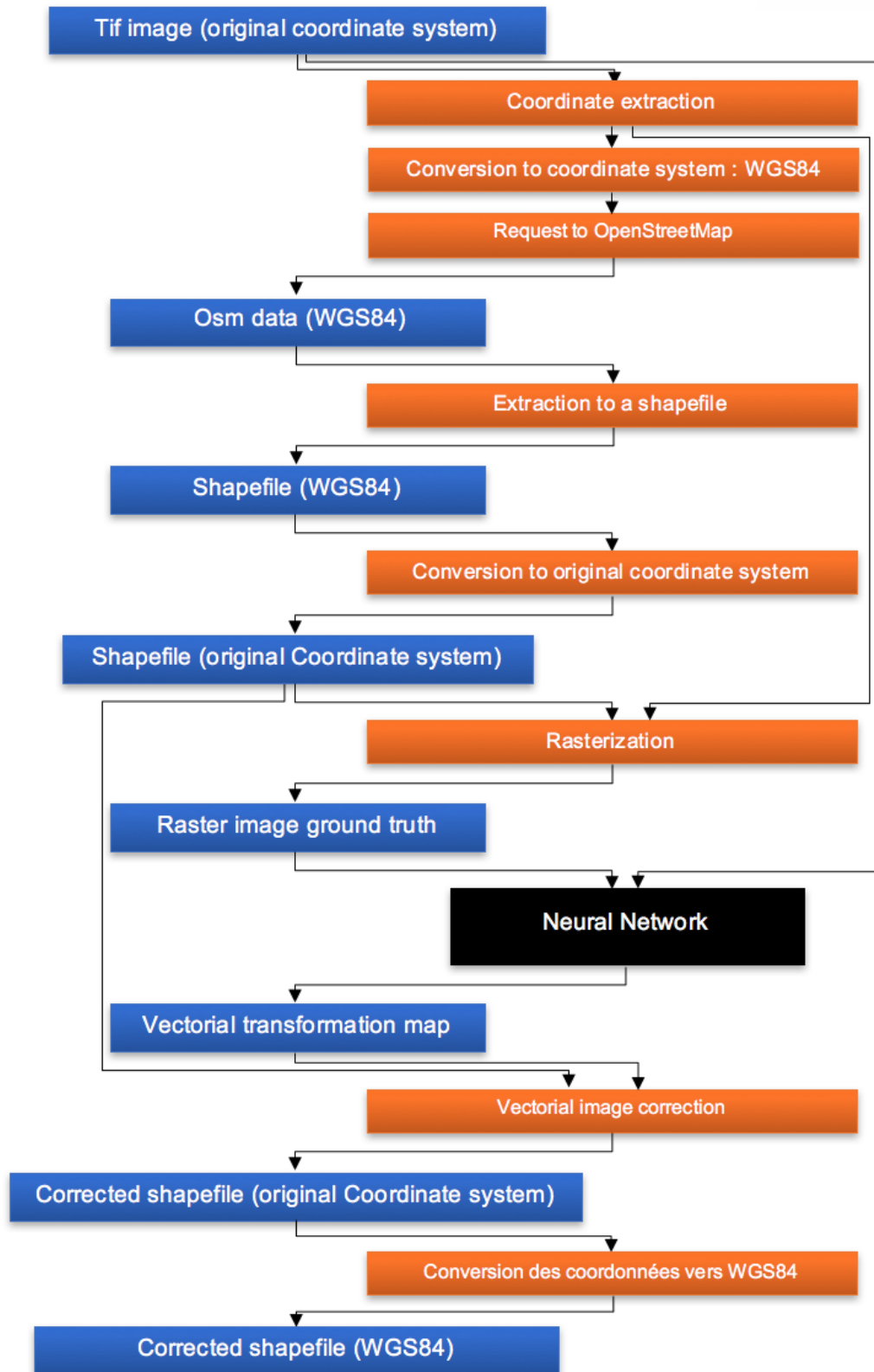


Figure 13. Global framework for OpenStreetMap data correction.

## A.4 Training details

### A.4.1 Tricks for better training

To reduce the training time and possible memory issues, patch of images ( $256 \times 256$  pixels) were given to the network for the training instead of the whole image, reducing the amount of computations needed per mini-batch. This is also important in terms of memory usage of the network as original images contain  $5000 \times 5000$  pixels. Furthermore, neural network computations and data generation (a random transformation is generated for each image at each training step in order to augment the training set size) are parallelized in order to improve the training speed of the algorithm.

Another issue encountered was to reach the local minimum corresponding to outputting always a null deformation, thus preventing the neural network parameters to evolve towards a better optimum. To solve this issue we used several methods to facilitate the training of the network and reduce its probability to reach this local minima. The first technique used was to force the network to overfit a very small sample of the dataset (400 iterations on 4 images with random transformations). This proved to be particularly efficient at avoiding the null local minimum.

The second technique is specific to the dataset used, *i.e.* the data from the cadastral images is particularly sparse. The first step is then to check, before selecting any training example, that data needed for alignment is present on the cadastral image, *i.e.* the cadastral image does contain buildings, *e.g.* This was done simply by calculating the ratio between labeled and unlabeled pixels (cadastral/non-cadastral) on each candidate image patch, and by setting a range of accepted values (*e.g.*, the minor class of an image should represent at least 5% of the labels of this image). Patches not respecting this rule were not selected when forming mini-batches. Thus mini-batches contained only relevant examples.

Another issue linked to the sparsity of the cadastre arises when the network is training on parts of a patch where not enough data is present to determine the transformation needed (*e.g.*, only one class within an area, as in a garden for example), even to a human eye: the estimation of the offset was not possible due to the lack of information. To solve this problem, we increase the weight of the loss function on the boundaries of buildings (parts where the transformation can be well estimated). For this, we first detect building boundaries based on the cadastre, as shown in 14, and add an multiplying factor to the loss at such locations (which is equivalent to sampling more often there). This insures that the deformation is findable and that the training is useful. This last trick is specific to our dataset, which is binary labelled, but we show in experiments that this is not needed

when the dataset is not sparse or binary.

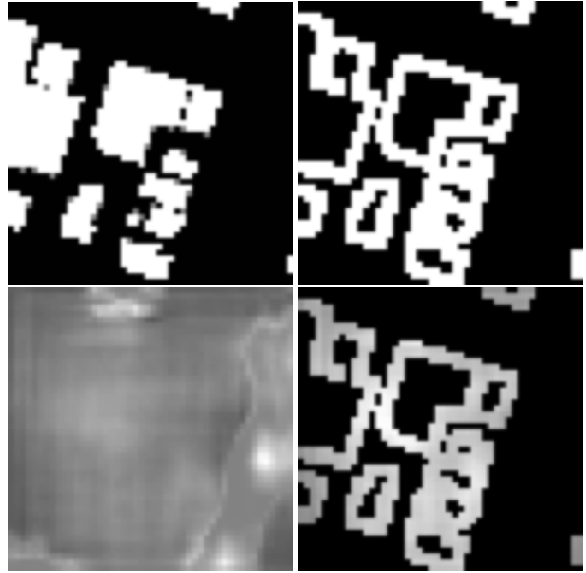


Figure 14. **Building boundaries used to re-weight the loss function.** Top left: cadastral image; top right: building boundary mask used; bottom left: the original map of the loss function; bottom right: masked map of the loss function on boundaries. The loss will be multiplied with a constant factor in these areas.

Lastly, we observed minor aligning issues when dealing with rows of houses with common walls, due to local translation invariance of the images, which adds locally a degree of freedom for alignment along the row axis. We decided to give supplementary information corresponding to the location of all corners extracted from the OpenStreetMap vectorial image (as each corner of each house is indicated), hoping to help to guide the alignment along such translation-invariant line in the cadastre. This step is however not critical as the results improvement is small and specific to certain building geometries (row of identical houses with common walls).

### A.4.2 Training information

Number of iterations : 60 000  
Batch size : 16  
Time to train : 16 hours  
Number of images : 108 original images (with a random transformation generated at each iteration for each image)  
Original image size :  $5000 \times 5000$   
Patch size :  $256 \times 256$   
Total number of layers : 26  
Memory used with tensorflow : 9.7 GB  
GPU : GeForce GTX 1080  
Processor : Dual-Xeon E5-2630  
RAM : 64 GB

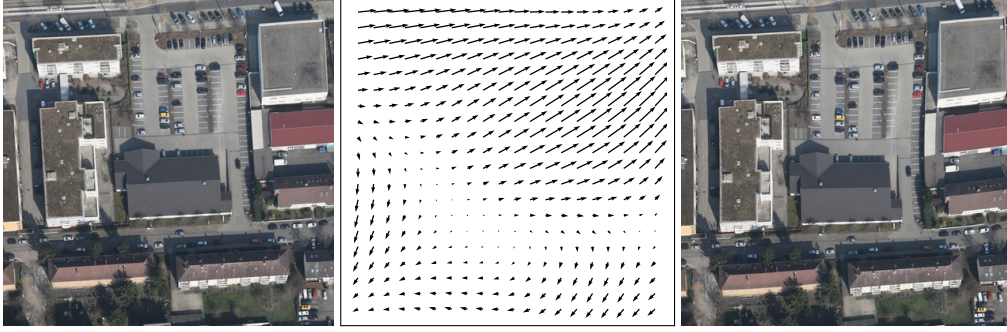


Figure 15. **Example of deformation.** Left: an image  $I$ ; middle: a deformation  $\phi$ , *i.e.* a  $\mathbb{R}^2$  vector field; right: the associated deformed image  $I \circ \phi$ .

## A.5 Keypoint matching

The keypoint matching experiment shown in the figure 7 of the main paper is shown full resolution here in Figure 16.

Supplementary matching examples are shown in Figure 17.

## A.6 Stereovision

In the paper we suggest to try the same framework on a different task, the one of stereovision. The result, shown in Figure 18 without any tuning, is very promising and confirms the generalization ability of our approach.

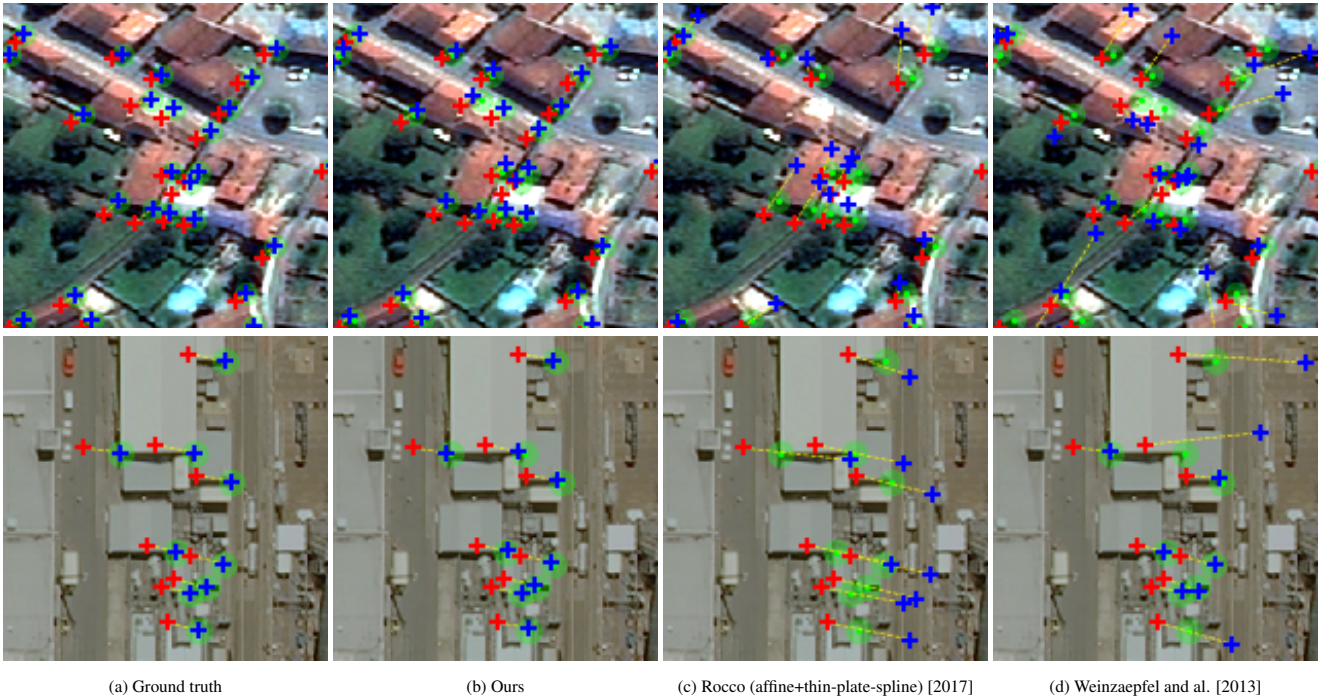


Figure 16. **Multimodal keypoint matching comparison** for different methods and two datasets. Top: Forez dataset; bottom: Kitsap. Blue: predicted, green: ground truth (centers of the green circles), red: original location of the corner (from the OpenStreetMap cadastral image, which is mis-geolocalized with respect to the RGB image).





Figure 17. **Additional multimodal keypoint matching examples.** Same setup as in Figure 16.





Figure 18. **Stereovision test.** Top: right image; middle: ground truth depth map for right image; down: predicted depth map.